

Edge detection:

image * filter = output

↑
convolution
operation

Vertical edge : filter

1	0	-1
1	0	-1
1	0	-1

take filter and do element wise product. move it along x + y axis

horizontal edge : filter

1	1	1
0	0	0
-1	-1	-1

Sobel filter :

1	0	-1
2	0	-2
1	0	-1

- ↳ i) more robust
- ii) more weight on center

Scharr filter :

3	0	-3
10	0	-10
3	0	-3

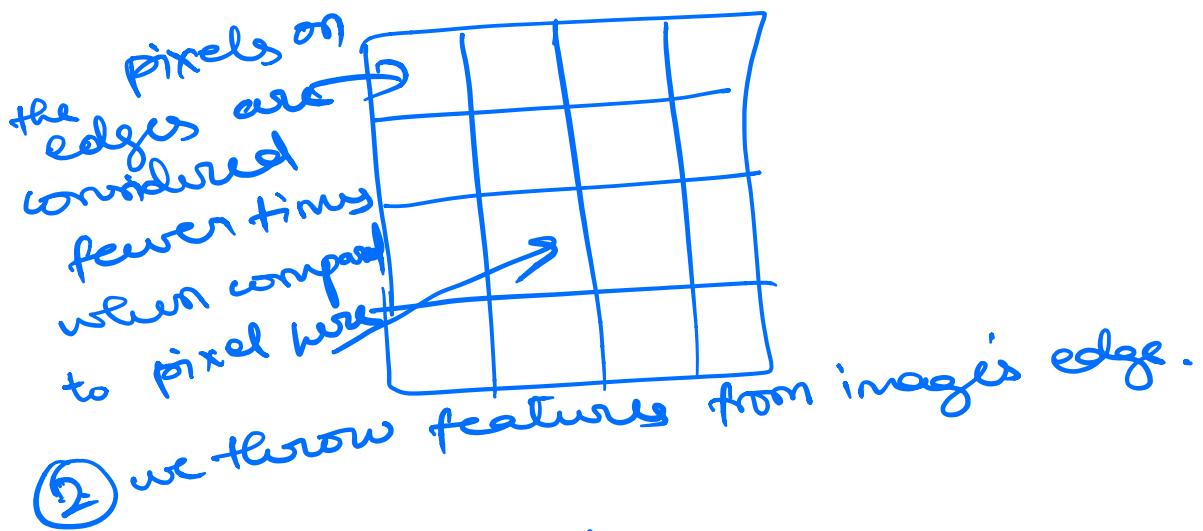
Padding :

input : $n \times n$

filter : $f \times f$

output : $n-f+1$,
 $n-f+1$

① Image shrinks



② we can pad.

$$6 \times 6 \rightarrow 8 \times 8$$

filter $3 \times 3 \rightarrow \text{output} = 6 \times 6$

out: $n+2P-f+1, n+2P-f+1$

\downarrow
padding size

valid convolution: no padding

$$n \times n * f \times f \rightarrow n-f+1 \times n-f+1$$

↑
convolve

same convolution: pad so that
output size
= input size

$$x + 2p - f + 1 = x$$

$$2p = f - 1$$

$p = \frac{f-1}{2}$ \Rightarrow ensures output
size same as input size

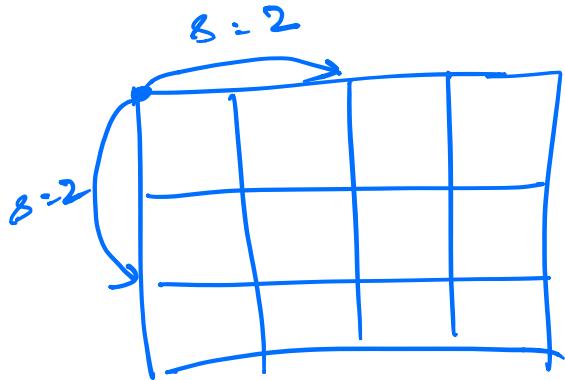
\therefore maintain size
of filter as odd

another reason to have f as odd

is: we have a center
for the filter as one point

Stride:

step by ' s ' steps instead of 1.



filter should fully
be inside image.
If not dont perform
operation.

$n \times n * f \times f$

padding P

stride s

$$\text{output} = \frac{n+2P-f}{s} + 1$$

round it down

Ureal convolution:

flip filter one time along
vertical and horizontal
axis.

If we dont flip, its called
cross-correlation

Convolutions over volume:

filter has third dimension

$h \times w \times \# \text{channel}$ $\xleftarrow{\text{input}}$

$h_f \times w_f \times \# \text{channel}$ $\xleftarrow{\text{filter}}$

the # channels in i/p + filter
should match.

$n \times n \times c * f \times f \times c \rightarrow \begin{matrix} n-f+1 \times \\ n-f+1 \times 1 \end{matrix}$

we can have multiple filters so:

$n \times n \times c * x (f \times f \times c) \rightarrow \begin{matrix} n-f+1 \times \\ n-f+1 \times x \end{matrix}$

\uparrow
no. of
filters

If layer l is a convolution layer:

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = number of filters

→ Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l]}$

Activations: $A^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias: $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$ #f: bias in layer l. $n_c^{[l]} \times n_H^{[l]} \times n_W^{[l]}$

Input: $\underbrace{n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}}_{n_{HW}^{[l-1]}} \leftarrow$
 Output: $\underbrace{n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}}_{n_{HW}^{[l]}} \leftarrow$

$$\frac{n_{HW}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1$$

$$A^{[l]} \rightarrow \underbrace{m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}}_{\text{number of neurons}}$$

Pooling layer:

parameters: f, s, p → usually zero

max pool: take max of input covered by filter

Average:
take average

It has no parameters to learn

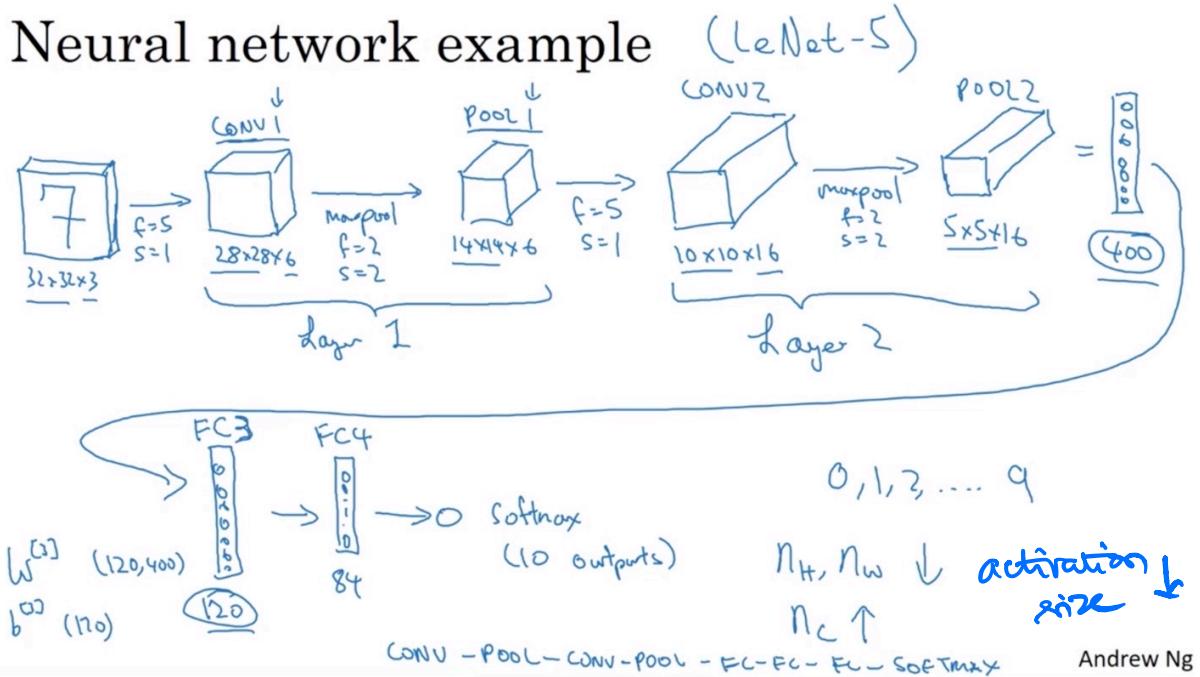
$$n, n, n_c \text{ maxpool } f, s \rightarrow \frac{n+2p-f+1}{s}$$

$$\frac{n+2p-f+1}{s}$$

for 3D:

compute max pool on each $\frac{n_c}{s}$ channels independently. same as input

layers = # layers with weight
 (so don't consider pool layer)

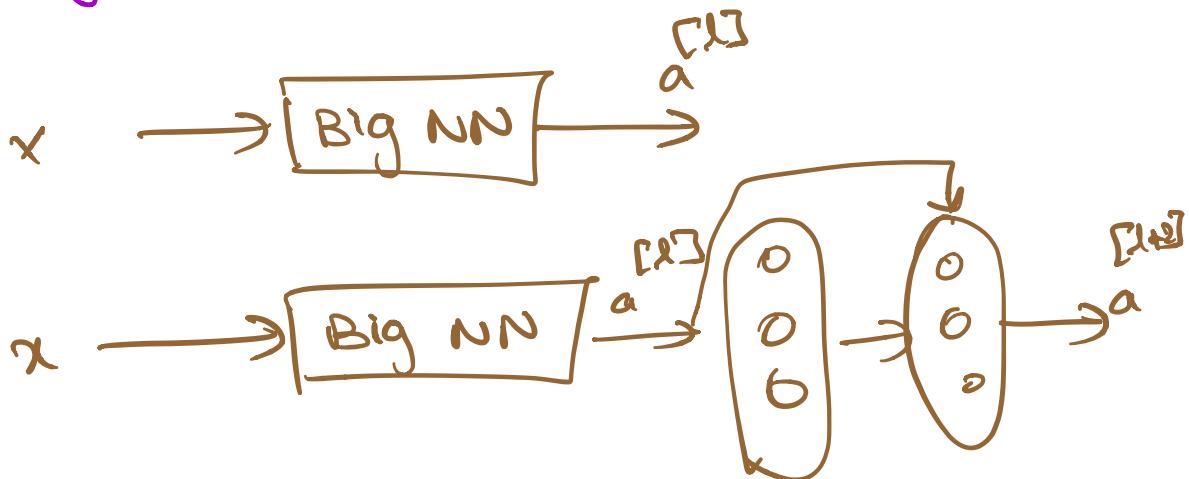


why convolution:
 → helps in parameter sharing
 → sparse connections
 (output value depends on small no. of inputs)

Advantages of having a pooling layer:

- 1) reduce size of representation
- 2) speed the computation
- 3) make detected features robust by maintaining the max.

Why residual networks - works?



we use value

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$$= g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]})$$

even if w and $b = 0$, then

$$a^{[l+2]} = g(a^{[l]}) = a^{[l]}$$

\therefore even adding 2 extra layers,

there won't be any issue/^{problem} in learning

as the last layer will learn identity fn. Therefore, it won't affect the network's performance.

And, as $a^{[l+2]}$ depends on $a^{[l]}$ and $a^{[l]}$, we will mostly use "SAME" padding in residual connections so that the shapes remains the same

If size not same

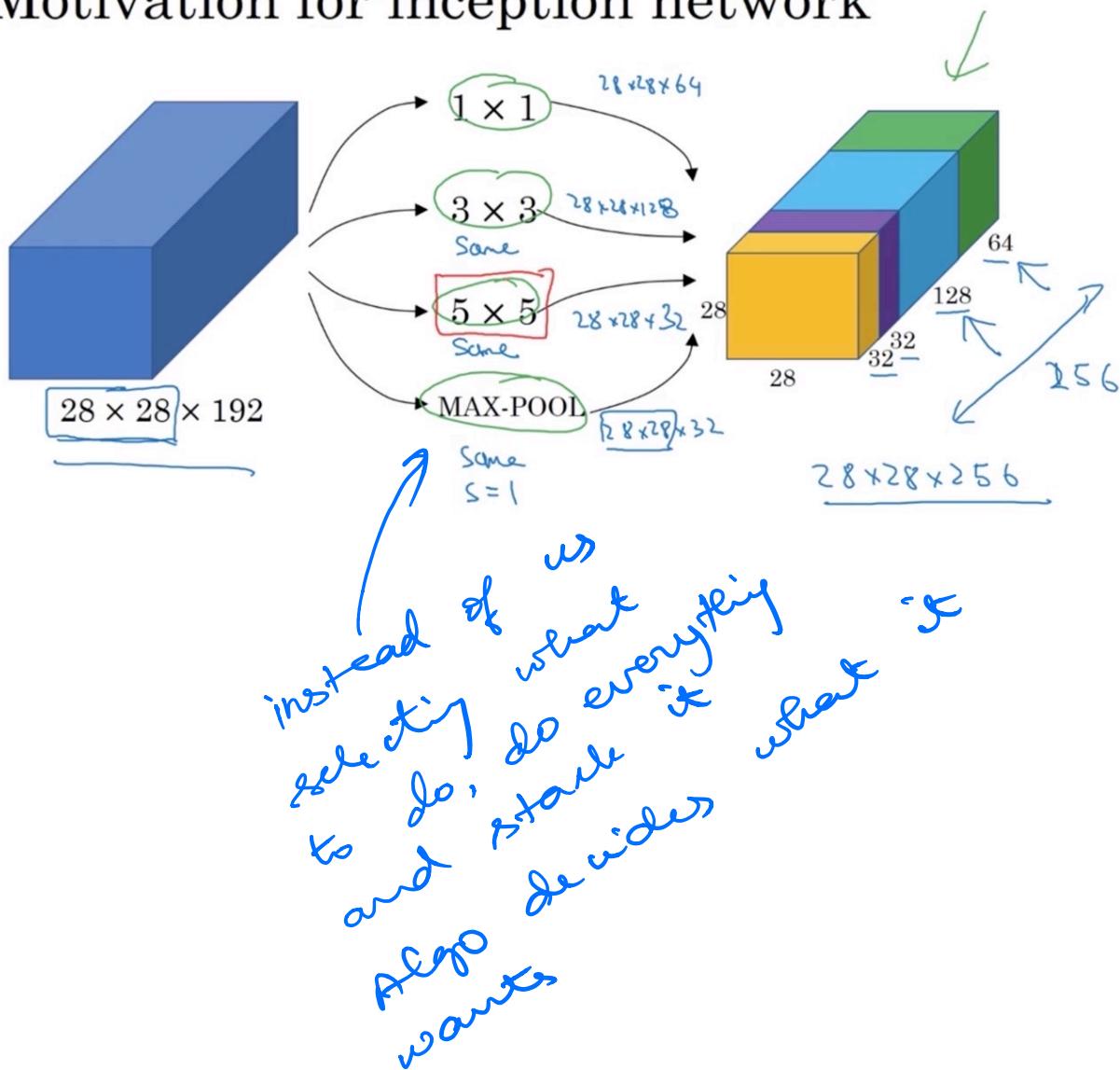
$$\Rightarrow a^{[l+2]} : g(z^{[l+2]} + w_3 a^{[l]})$$

learned or fixed (zeroed)

Network in network (1x1 convs):

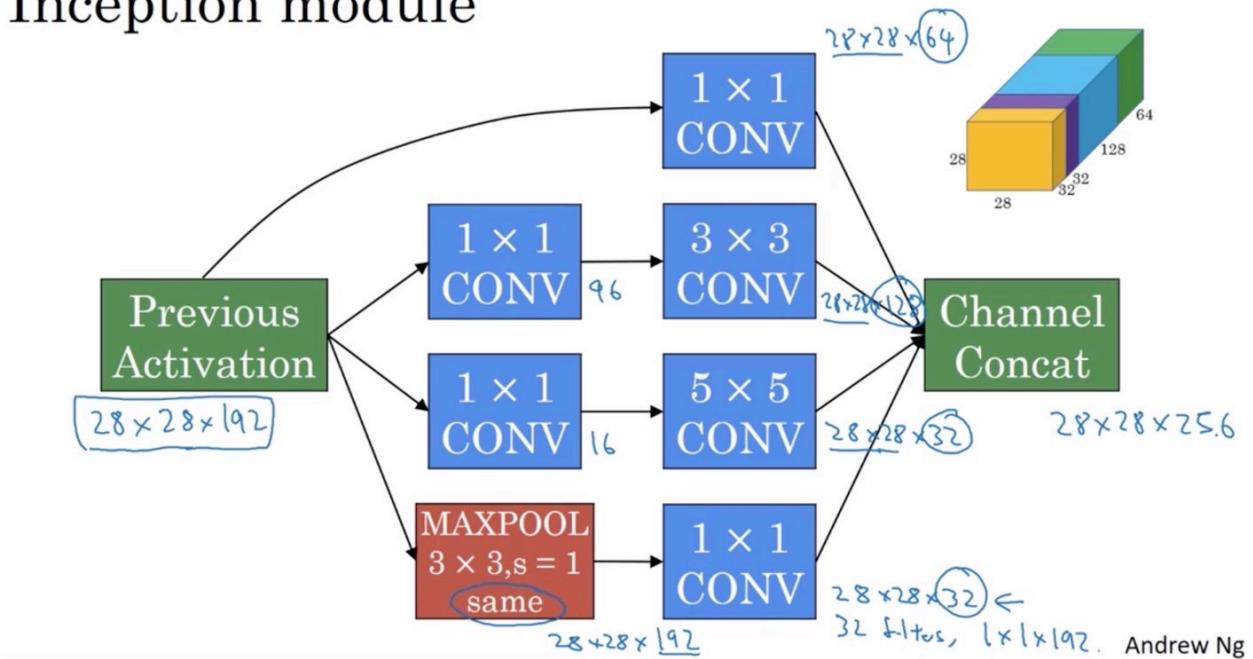
acts as a single neuron
that takes all inputs \times weights,
and applies relu.

Motivation for inception network



To overcome computation cost, we a 1×1 conv to reduce the input size and apply other conv layer on reduced input. This decreases # comput by a large amount.

Inception module



Data augmentation :

- i) flip image iii) random cropping
- ii) shearing iv) color shifting

Localization:

Tell where the object is

detection:

multiple objects of different category

Localization:

make the conv net output

4 more values. ($\underbrace{bx, by, bh, bw}_{\substack{\text{midpt} \\ \text{of bound} \downarrow \\ \text{height} \downarrow \\ \text{width} \downarrow}})$

probability that
an object (not background)
is present

$$y =$$

there are
g elements

P_c
bx
by
bh
bw
c_1
c_2
c_3

is there any
object (other than
background)

classes in the
image

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots \\ \quad + (\hat{y}_B - y_B)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

In practice, we can use:

- i) log like feature loss for c_1, \dots, c_n
to the softmax output
- ii) squared error for bounding box
- iii) logistic regression loss for p_c

we can also make a CNN to output landmark / exact x,y point in an image

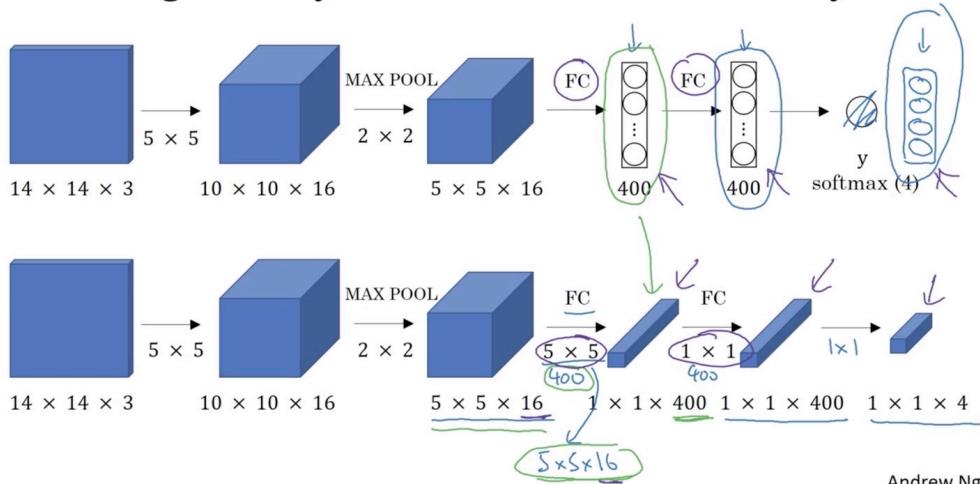
Object detection:

- i) train a convnet on closely cropped images of the object.
- ii) Use sliding windows on the input image and detect object in it using above convnet.
→ do this repeatedly for varying window size

This is called sliding window detection

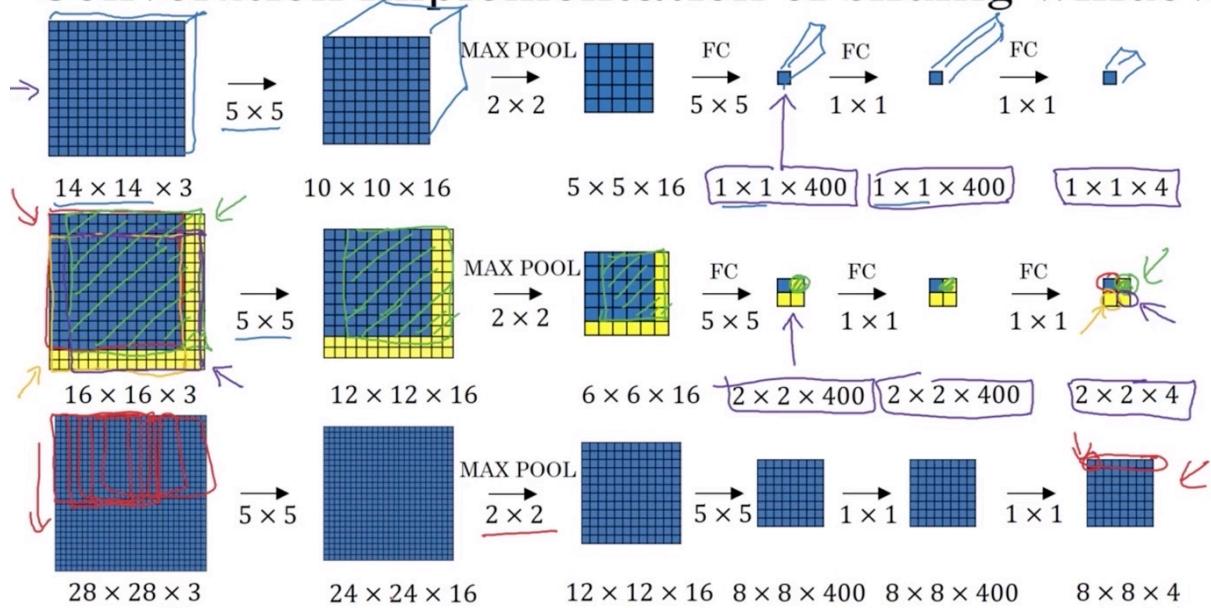
Computational cost is very high

Turning FC layer into convolutional layers



Andrew Ng

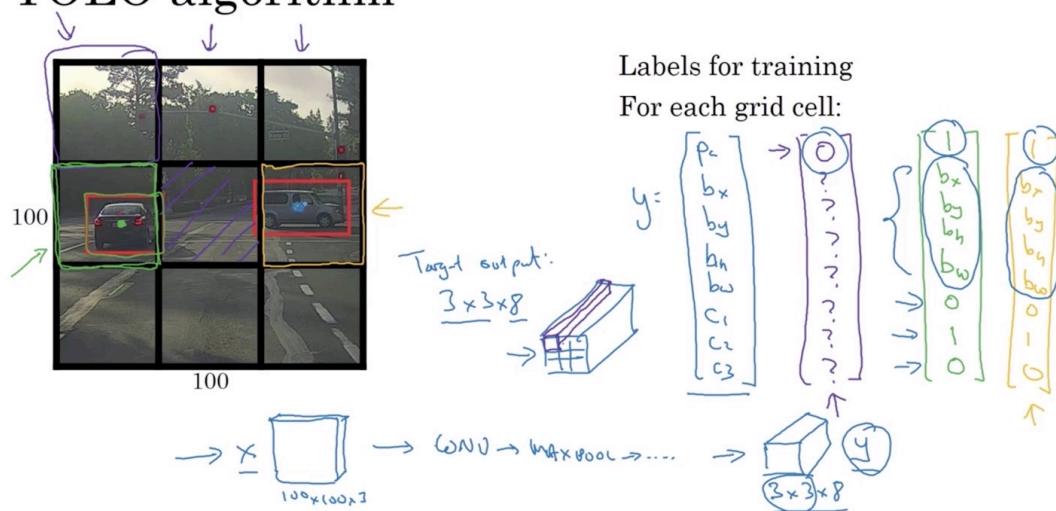
Convolution implementation of sliding windows



Still the predictions are not ~~so~~
accurate!

You only look once algo: (YOLO)

YOLO algorithm



Summary

- i) split image into grids
- ii) use object localization in each grid.

This can be done using a fully conv net and output $n \times n \times 8$ ← corresponds to 8 value in obj. loc. output.
↑
no. of grids

here in Yolo, each grid's starting point is $(0,0)$ and end is $(1,1)$
so b_x, b_y, b_h, b_w are relative to grid cell

Intersection over union:

$$= \frac{\text{size of intersection}}{\text{size of union}}$$

"correct" if $\text{IoU} \geq 0.5$

Non max suppression:

Get $19 \times 19 \times 8$ output after running object detection.

Consider each bounding box has a probability / confidence. Discard $P_c \leq 0.6$

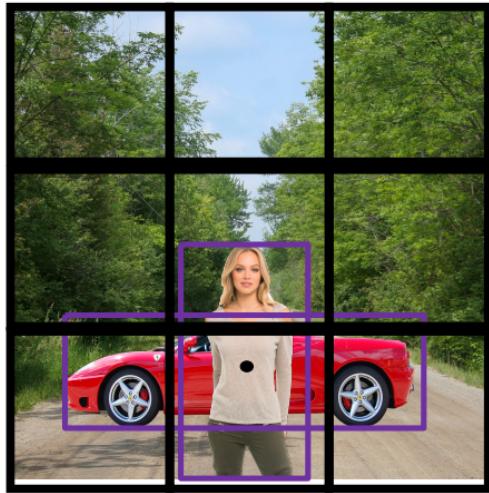
Take the max prob.

Find other boxes that has high iou (≥ 0.5) with this box and suppress it.

Continue until we have non overlap box. (i.e) until we have some remaining boxes.

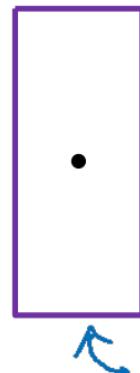
If we have multiple objects, do non-max suppression independently on each box.

Overlapping objects:



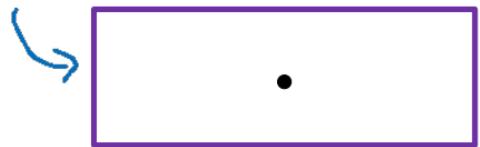
$$y = \begin{bmatrix} P_C \\ b_x \\ b_y \\ \vdots \\ c_3 \end{bmatrix}$$

Anchor box 1:



$y =$

Anchor box 2:



Anchor box 1

Anchor box 2

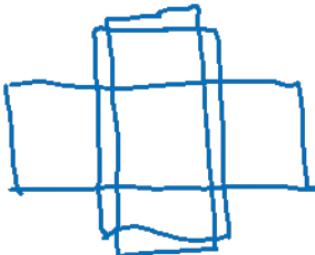
Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output y:

$3 \times 3 \times 8$



With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

(grid cell, anchor box)

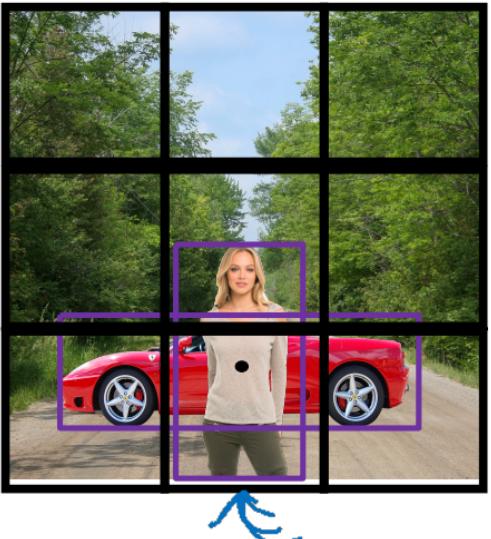
Output y:

$3 \times 3 \times 16$

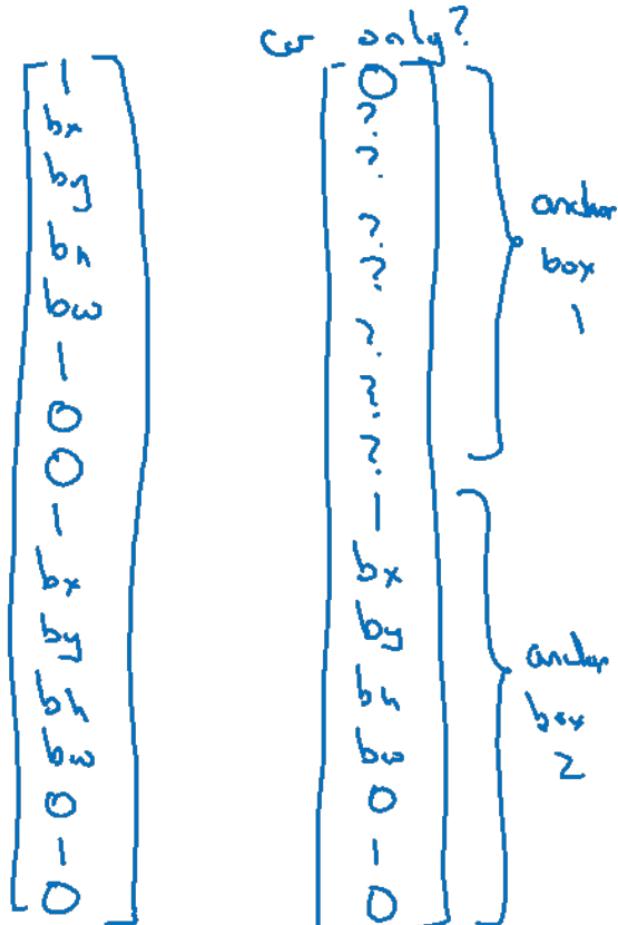
$3 \times 3 \times 2 \times 8$

Anchor box example

y =



Anchor box 1: Anchor box 2:



RCNN:

instead of using sliding window on all possible window, select regions on which to run convnet classifier.

First run segmentation algo to get segmented image and we convnet on a blob in segmented image to get label and bounding box.

This method is slow

Fast-R-CNN : we convolutional sliding window.

Faster-R-CNN: use CNN for segmentation

face recognition:

face verification

input image, name/ID
output: if input image
is that of claimed
person

↓
1:1

vs recognition

↓

- * has a dataset of K person
- * input image
- * output ID if image is any of K person

↓
1:K

one shot learning:

learn from one example to
recognize the person again

general approach: (4 persons)



but what if new person joins?
→ we have to train whole
network with 6 op units

so, learn similarity fn:

$d(\text{img}_1, \text{img}_2)$: degree of diff.
between images

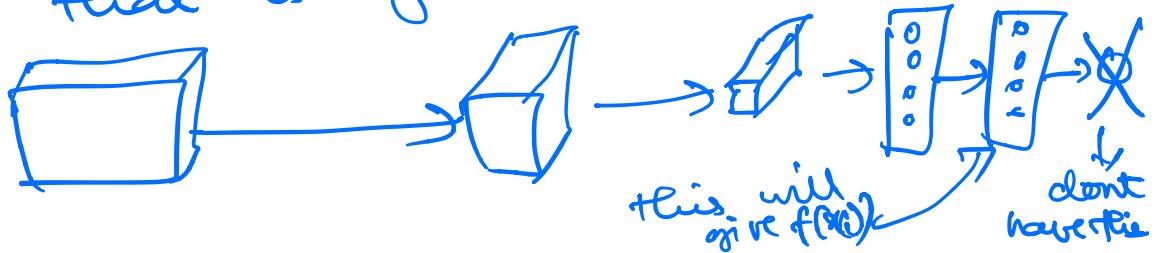
during test:

if $d(\text{img}_1, \text{img}_2) \leq T \rightarrow \text{same person}$
 $> T \rightarrow \text{diff person}$

compute $d(\text{img}_1, \text{img}_2)$ if $\text{img}_2 \in \text{dataset}$
and see id where it's
small.

Siamese network:

train a NN with last layer
giving the encoding of image
that is given as input.



learn parameters such that:

If $x^{(i)}, x^{(j)}$ are same:
 $\|f(x^{(i)}) - f(x^{(j)})\|_2^2$
is small

if $x^{(i)}, x^{(j)}$ are different:

$\|f(x^{(i)}) - f(x^{(j)})\|_2^2$
is large

Triplet loss:

we have an anchor image
and a positive image.
 \Rightarrow distance has to be small

anchor image and a negative
image
 \Rightarrow distance has to be large

Triplet loss: look 3 images
Anchor, positive, negative
A P N

$$\text{want : } \left\| f(A) - f(P) \right\|_2^2 \\ \leq \left\| f(A) - f(N) \right\|_2^2$$

$$d(A, P) \leq d(A, N)$$

$$\left\| f(A) - f(P) \right\|_2^2 - \left\| f(A) - f(N) \right\|_2^2 \leq 0$$

$$\left\| f(A) - f(P) \right\|_2^2 - \left\| f(A) - f(N) \right\|_2^2 + \alpha \leq 0$$

this margin is added
so that our CNN doesn't
output 0 for all encoding

loss fn:

given : A, P, N (3 images)

$$L(A, P, N) = \max \left(\left\| f(A) - f(P) \right\|_2^2 + \left\| f(A) - f(N) \right\|_2^2 + \alpha \right)$$

we take max because
term A should be ≤ 0 .
If it's ≤ 0 , loss is 0

Since we need A,P to be same person, for training we need multiple pictures of same person. But while testing, we need only one image per person.

If A,P,N are chosen randomly (A:P)
then it will be easy to satisfy triplet loss fn and our NN wont learn anything.

∴ we need to choose triplets that's hard to train on

$$d(A, P) + \alpha \leq d(A, N)$$

$$d(A, P) \gtrsim d(A, N)$$

⇒ learning becomes hard

we can also treat this as a binary classification where we feed $f(x^{(i)})$ and $f(x^{(j)})$ to a logistic regression and output 1 (same) or 0 (diff. person)

$$y = \sigma \left(\sum_{k=1}^{128} w_k \underbrace{|f(x^{(i)})_k - f(x^{(j)})_k|}_{+ b} \right)$$

$$\chi^2 \text{ (chi square)} = \frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$$

Neural style transfer:

content + style \rightarrow generated
 (C) (S) $\overset{\text{image}}{(G)}$

$$J(G) = \alpha J_{\text{content}}(C, G)$$

find G : $+ \beta J_{\text{style}}(S, G)$

- i) randomly initiate G
- ii) use gradient descent to minimize $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G)$$

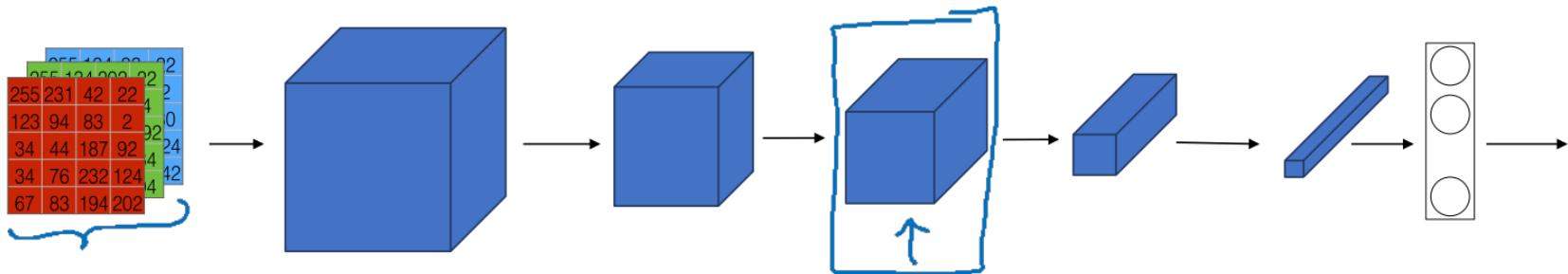
Content cost function

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

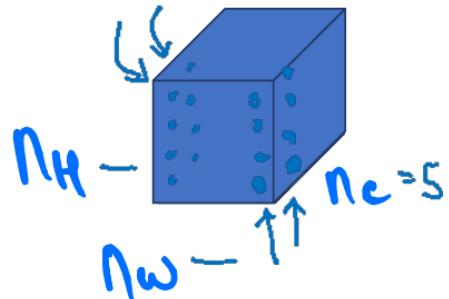
- Say you use hidden layer l to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let $a^{[l](C)}$ and $a^{[l](G)}$ be the activation of layer l on the images
- If $a^{[l](C)}$ and $a^{[l](G)}$ are similar, both images have similar content

$$J_{\text{content}}(C, G) = \frac{1}{2} \| \underbrace{a^{[l](C)}}_{\uparrow} - \underbrace{a^{[l](G)}}_{\downarrow} \|_2^2$$

Meaning of the “style” of an image



Say you are using layer's activation to measure "style."
Define style as correlation between activations across channels.

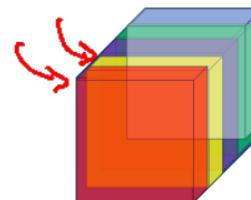


How correlated are the activations
across different channels?

Intuition about style of an image



Generated Image



Correlated?
Un correlated → Components appears together
← So in generate image's activation, we can find correlated styles and use it.

Style matrix:

Let $a_{i,j,k}^{[l]} = \text{activation at } (i,j,k)$

$G_l^{[l]}$ is $n_c \times n_c$

↑
each pair of channels' correlation

gram matrix $\rightarrow G_l^{[l]}$ → correlation of channel k and k'

$$G_{kk'}^{[l]} = \sum_i^H \sum_j^W a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

iterate over
o to n_c

DO $G_{kk'}^{[l]}$ for both style and generated image

$$\mathcal{J}_{\text{Style}}(S, G) := \frac{1}{n_w n_h n_c} \| G^{[l](S)} - G^{[l](G)} \|_F^2$$

$$= \frac{1}{(2 n_w n_h n_c) K} \sum_{k'} \left(G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)} \right)^2$$

$$J_{\text{style}}(S, G) = \sum_l J_{\text{style}}^{[l]}(S, G)$$

↑

iterate over
diff layers for
getting pleasing result